



**University of
Zurich**^{UZH}

**Zurich Open Repository and
Archive**

University of Zurich
University Library
Strickhofstrasse 39
CH-8057 Zurich
www.zora.uzh.ch

Year: 2021

Why we learn less from observing outgroups

Kang, Pyungwon ; Burke, Christopher John ; Tobler, Philippe N ; Hein, Grit

Abstract: Humans are less likely to learn from individuals belonging to a different group (outgroup) than from individuals of their own group (ingroup), yet the source of this societally relevant deficit has remained unclear. Here we used neuroimaging and computational modeling to investigate how people learn from observing the actions and outcomes of ingroup and outgroup demonstrators. Politically left-wing male and female participants performed worse when observing computer-simulated actions they believed were from a right-wing outgroup member compared with those from a left-wing ingroup member. A control experiment in which participants observed choices from a nonhuman agent confirmed that this performance difference reflected an outgroup deficit, rather than an ingroup gain. Accounting for the outgroup deficit, a computational model showed that participants relied less on information from outgroup actions compared with ingroup actions, while learning from outgroup outcomes was not impaired. At the neural level, the differences in observational ingroup versus outgroup learning were reflected in lateral prefrontal activity. The stronger the activity in this region, the more strongly participants weighed ingroup compared with outgroup learning signals (action prediction errors), which formally captured deficits in outgroup learning. Together, our work provides a computational and neural account of why people learn less from observing outgroups.

DOI: <https://doi.org/10.1523/jneurosci.0926-20.2020>

Posted at the Zurich Open Repository and Archive, University of Zurich

ZORA URL: <https://doi.org/10.5167/uzh-193175>

Journal Article

Accepted Version



The following work is licensed under a Creative Commons: Attribution 4.0 International (CC BY 4.0) License.

Originally published at:

Kang, Pyungwon; Burke, Christopher John; Tobler, Philippe N; Hein, Grit (2021). Why we learn less from observing outgroups. *Journal of Neuroscience*, 41(1):144-152.

DOI: <https://doi.org/10.1523/jneurosci.0926-20.2020>

Research Articles: Behavioral/Cognitive

Why We Learn Less from Observing Outgroups

<https://doi.org/10.1523/JNEUROSCI.0926-20.2020>

Cite as: J. Neurosci 2020; 10.1523/JNEUROSCI.0926-20.2020

Received: 14 April 2020

Revised: 2 November 2020

Accepted: 4 November 2020

This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.

Alerts: Sign up at www.jneurosci.org/alerts to receive customized email alerts when the fully formatted version of this article is published.

Why we learn less from observing outgroups

Running title: Outgroup deficit in observational learning

Pyungwon Kang¹, Christopher J. Burke¹, Philippe N. Tobler^{1*}, & Grit Hein^{2*}

¹ Department of Economics and Laboratory for Social and Neural Systems Research, University of Zurich

² Translational Social Neuroscience Unit, Department of Psychiatry, Psychosomatic and Psychotherapy, University of Würzburg

* shared senior authorship

Correspondence to: pyungwon.kang@gmail.com, Blumlisalpstrasse 10, 8006, Zurich, Switzerland

Number of pages: 34

Number of figures: 6

Number of tables: 2

Number of words for abstract: 203

Number of words for introduction: 642

Number of words for discussion: 1395

Conflict of interest statement: The authors declare no competing financial interests.

Acknowledgements: This work was supported by the Swiss NSF (IZKSZ3_162109 and 100019_176016 to PNT). We also acknowledge funding by the German Research Foundation (HE 4566/5-1; HE 4566/3-1 to GH).

29 **Abstract**

30 Humans are less likely to learn from individuals belonging to a different group (outgroup) than
31 from individuals of their own group (ingroup), yet the source of this societally relevant deficit
32 has remained unclear. Here we used neuroimaging and computational modeling to investigate
33 how people learn from observing the actions and outcomes of ingroup and outgroup
34 demonstrators. Politically left-wing male and female participants performed worse when
35 observing computer-simulated actions they believed were from a right-wing outgroup member
36 compared to those from a left-wing ingroup member. A control experiment in which participants
37 observed choices from a non-human agent confirmed that this performance difference reflected
38 an outgroup deficit, rather than an ingroup gain. Accounting for the outgroup deficit, a
39 computational model showed that participants relied less on information from outgroup actions
40 compared to ingroup actions, while learning from outgroup outcomes was not impaired. At the
41 neural level, the differences in observational ingroup vs outgroup learning were reflected in
42 lateral prefrontal activity. The stronger the activity in this region, the more strongly participants
43 weighed ingroup compared to outgroup learning signals (action prediction errors), which
44 formally captured deficits in outgroup learning. Together, our work provides a computational
45 and neural account of why people learn less from observing outgroups.

46

47

48

49

50

51

52 **Significance statement**

53 Learning from observing others is an efficient way to acquire knowledge. In our globalized
54 world, “the others” often are people from a different social group (outgroup). There is evidence
55 that people learn less from observing outgroup individuals compared to individuals from their
56 own group (ingroup). However, the source of this outgroup deficit in observational learning
57 remained unknown, which limits our chances to improve intergroup learning. Our results showed
58 that participants rely less on observed outgroup actions compared to ingroup actions, while
59 learning from outgroup outcomes is not impaired. On the neural level, this outgroup deficit was
60 reflected in activation of the inferior frontal gyrus. These findings imply that intergroup learning
61 should rely on observing outcomes, rather than actions.

62

63

64

65

66

67

68

69

70

71

72

73

74

75 **Introduction**

76 For many different species, including humans, learning from perceiving the actions and
 77 outcomes of others (i.e., observational learning) is an efficient way to acquire knowledge and
 78 skills. There is evidence that observational learning is modulated by important social factors such
 79 as group membership (Buttelmann et al., 2013; Golkar et al., 2015; Howard et al., 2015). For
 80 example, learning appears to be facilitated if participants observe a person from their own social
 81 group (ingroup) compared to a person from a different social group (outgroup) (Golkar et al.,
 82 2015; Golkar and Olsson, 2017). This ingroup bias in observational learning was even found in
 83 infants and children being more likely to imitate the novel actions of a demonstrator who speaks
 84 their language, compared with a person speaking a different language (Buttelmann et al., 2013;
 85 Howard et al., 2015).

86 Previous neuroscience studies have investigated observational learning irrespective of
 87 group membership (Burke et al., 2010; Suzuki et al., 2012; Charpentier et al., 2020; Kumaran et
 88 al., 2015) and reported learning from observed outcomes and from observed actions. Learning
 89 from others' outcomes and the resulting outcome prediction errors is associated with activation
 90 in medial prefrontal cortex (MPFC; Burke et al., 2010; Suzuki et al., 2012; Kumaran et al., 2015).
 91 Learning from others' actions and the resulting action prediction errors is related to activation in
 92 dorsolateral prefrontal cortex/inferior frontal gyrus (DLPFC/IFG; (Burke et al., 2010; Suzuki et
 93 al., 2012; Charpentier et al., 2020).

94 The effect of group membership has mainly been investigated with regard to action
 95 observation or imitation (Losin et al., 2012; Losin and Woo, 2015), revealing ingroup vs
 96 outgroup differences in brain regions associated with mentalizing (Losin and Woo, 2015), and
 97 parts of the 'mirror neuron system' such as the inferior frontal, motor and parietal cortex (Losin

et al., 2012). However, it remained unclear how important social factors, such as group membership, shape observational learning mechanisms and the underlying neural circuitries.

In our study, we investigated how group membership affects the neural circuitries of observational learning. To do so, we combined a well-established observational learning paradigm (Burke et al., 2010) with social group manipulation, computational modeling and fMRI. In more detail, participants inside the fMRI scanner observed only choices (i.e., actions) or choices and outcomes of an ingroup and an outgroup demonstrator, and could use these different pieces of information to optimize their own choice.

Based on previous behavioral evidence of outgroup deficits in social learning (Buttelmann et al., 2013; Golkar and Olsson, 2017), we hypothesized that participants choose the “correct” (i.e., more rewarding) option less frequently after observing outgroup choices compared to ingroup choices, reflecting an outgroup deficit in observational learning. Given that individuals learn from observing A) the outcomes, and B) the actions of others (Burke et al., 2010; Suzuki et al., 2012), we derived three different hypotheses regarding the mechanisms that might underlie the potential deficit in outgroup learning. According to a first hypothesis, the outgroup deficit in observational learning may arise because participants rely more on observed ingroup compared to outgroup outcomes. In computational modeling, this should be reflected by a stronger weight for ingroup compared to outgroup outcome prediction errors, associated with neural activation of the MPFC (Burke et al., 2010; Suzuki et al., 2012; Kumaran et al., 2015). Alternatively, the outgroup deficit in observational learning may occur because participants rely more on observed ingroup compared to outgroup actions. In this case, our computational modeling results should reveal a stronger weight for ingroup compared to outgroup action prediction errors, related to increased activation in DLPFC/IFG (Burke et al., 2010; Suzuki et al.,

121 2012; Charpentier et al., 2020). Finally, it is possible that the outgroup deficit in observational
 122 learning arises because participants rely more on observing ingroup outcomes and actions,
 123 reflected by a stronger weight for ingroup compared to outgroup outcome and action prediction
 124 errors, paralleled by increased activation in MPFC and DLPFC/IFG.

125

126 **Methods**

127 **Participants**

128 fMRI study. Thirty-two participants (19 female, mean age: 22.51 ± 0.54 years) were recruited
 129 from the University of Zurich and a local community in Zurich. Participants were all right-
 130 handed, had normal/corrected-to-normal vision and did not have a history of psychological or
 131 neurological disorder. Because our group manipulation was based on political attitude (see
 132 below), we invited Swiss participants who perceived themselves as politically active (i.e.,
 133 interested in current political debates in Switzerland) with primarily left-wing attitudes. Three
 134 participants had to be excluded because they showed right-wing attitudes during the group
 135 manipulation check (see below). Thus, we analyzed data of 29 participants in the imaging
 136 experiment.

137 Behavioral control study. For a behavioral control experiment outside the fMRI scanner,
 138 we recruited another sample of 33 participants (female: 18, mean age: 23.25 ± 3.34). The
 139 participants were recruited from the same participant pool and they matched the participants of
 140 the fMRI study in terms of age, education level, political attitude, and nationality, all p s > 0.171 .
 141 All participants received a fixed monetary compensation for their participation and additional
 142 incentives according to their performance. The study was approved by the ethics committee of
 143 the Canton of Zurich.

144

145 **Experimental design and statistical analysis**

146 fMRI study. Prescanning procedure for group induction. Before the main fMRI experiment,
 147 participants provided their political views by rating current political issues in Switzerland, and
 148 observed the ostensible ratings of two other individuals. One of these individuals displayed
 149 similar ratings as the participant, indicating a left-wing attitude. The ratings of the other
 150 individual indicated attitudes opposite to those of the participants (i.e., right-wing attitudes).

151 The prescanning procedure consisted of six trials. In each of these trials the participants
 152 were presented with a political initiative dividing the left- and right-wing parties (e.g., an
 153 initiative to raise inheritance tax). They were asked to indicate their opinion by moving an
 154 abstract symbol on a visual analogous scale (ranging from “strongly disagree” to “strongly
 155 agree”; **Figure 1A**). Next, participants were presented with two abstract symbols on the same
 156 rating scale that ostensibly indicated the ratings of two different individuals (**Figure 1A**). One of
 157 these symbols (designated to become the symbol of the ingroup demonstrator) appeared in the
 158 part of the rating scales that indicated agreement with left-wing initiatives and disagreement with
 159 right-wing initiatives, corresponding to a political attitude similar to that of the participants. The
 160 other symbol (designated to become the symbol of the outgroup demonstrator) appeared in the
 161 part of the rating scales that indicated agreement with right-wing initiatives and disagreement
 162 with left-wing initiatives, opposing the political attitude of the participant. The presentation order
 163 of the symbols was randomized across trials. The symbols representing the ingroup and the
 164 outgroup demonstrator were counterbalanced across participants, but remained constant within
 165 each participant. At the end of the group induction, the ratings of participants and the two other
 166 individuals (i.e., all three symbols) were again shown for each political issue to remind the

167 participants of the political attitudes of the other two individuals relative to their own attitude
 168 (**Figure 1A**).

169 Next, participants were asked to rate how close they feel to major political parties in
 170 Switzerland ranging from left-wing to right-wing, and to provide the same closeness ratings for
 171 each of the two other individuals whose ratings regarding political initiatives they had observed
 172 before. To do so, participants moved the respective symbols on a rating scale (ranging from
 173 “very close” to “not close at all”). These ratings served as manipulation check of our group
 174 manipulation because they quantified how differently the participants perceived the person
 175 associated with the ingroup and the outgroup symbol, and verified that the perceived differences
 176 resulted in social categorization (here supporters of left-wing and right-wing political parties).

177 *Observational learning task.* We used a modified version of an observational learning
 178 task established in a previous study (Burke et al., 2010). The participants were instructed to learn
 179 about the reward probability of two fractal images through observation of only choices
 180 (action_only condition) or choices and outcomes (action_outcome condition) of the two
 181 demonstrators whose political attitudes they had rated outside the scanner (**Figure 1A**). To do so,
 182 they would observe prerecorded choices that the two individuals (demonstrators) made to
 183 optimize their own outcomes in a previous experiment. Unknown to the participant, all
 184 demonstrator choices were generated by a standard reinforcement learning algorithm with a
 185 learning rate of 0.3 and a beta of 0.4 in each condition. We used this manipulation to ensure that
 186 participants observed identical choices from the ingroup and the outgroup demonstrator, action-
 187 outcome condition, $F(1,172) = 0.101$, $p = 0.751$, $\eta^2 = 0.001$ and in the action only condition, $F(1,$
 188 $172) = 0.201$, $p = 0.655$, $\eta^2 = 0.001$. Moreover, estimating the actually realized betas of ingroup
 189 and outgroup demonstrators revealed no difference between conditions, $t(28) = .91$, $p = 0.36$.

190 Thus, participants observed comparable choices (and a comparable learning rate) in the ingroup
 191 and the outgroup conditions, and potential differences between the conditions should result from
 192 the group manipulation (i.e., the differences in the demonstrators' group membership).

193 Each trial of the observational learning task consisted of an observation phase (i.e.
 194 observing the demonstrator's decision) and a decision phase (i.e. making a decision for
 195 themselves; **Figure 1C**). To prevent confusion between the two phases, the screen was vertically
 196 split in two halves with one half showing the observation phase and the other half showing the
 197 decision phase. The display side of the observation and decision phase was constant within each
 198 participant, and counterbalanced across participants.

199 At the beginning of the observation phase, the demonstrator for the present trial was
 200 indicated by one of the two previously learned abstract symbols, presented for a jittered duration
 201 of 1-10s. Then, participants had 1.5 s to predict which option the demonstrator would choose.
 202 After response registration (0.5s) and a short interval (0.5s), the action of the demonstrator was
 203 shown for 1s. Depending on the condition, participants then observed the outcome of the
 204 demonstrator (action_outcome condition) or a pixel-matched scrambled image (action_only
 205 condition) for 1.5 s. The scrambled image was made from the outcome presentation part (0 or 10)
 206 and all pixels of the respective outcome were randomly redistributed. After a jittered interval of
 207 1-10s, the demonstrator's symbol was replaced by the participant's symbol, followed by the two
 208 fractals (~1.5s). Within this period, the participants selected one of the two fractals. The choice
 209 was displayed for 0.5s, followed by a scrambled outcome image (1.5s), which prevented
 210 individual learning during the task. Learning was incentivized by paying out participant
 211 decisions at the end of the experiment. At the end of the experiment, participants were debriefed
 212 and filled in a questionnaire asking open questions regarding the purpose of the study and how

213 they felt during the study in the fMRI scanner. None of the participants reported suspicions about
 214 the experimental setting or inferred the purpose of the study (effect of group membership on
 215 observational learning) correctly.

216 The fractals were associated with different reward contingencies (70% reward vs. 30%
 217 reward). We used the two conditions (action-outcome vs. action-only) to disentangle action-
 218 related learning and outcome-related learning. In the action-outcome condition, participants
 219 could learn from both actions and outcomes of the demonstrators, while in the action-only
 220 condition they could learn only from observing the actions of the demonstrators. The
 221 observational learning task consisted of three ingroup condition and three outgroup condition
 222 sessions, with each session comprising one action_outcome block and one action_only block.
 223 The sequence of ingroup and outgroup sessions was interleaved, and the order of sessions was
 224 counterbalanced across participants.

225 *Individual learning task.* In addition to the observational learning task, participants
 226 performed an individual learning task. The structure of this task was identical to the
 227 observational learning task. However, now the symbol representing the demonstrator was
 228 replaced by a scrambled image, and participants were asked to press a random key instead of
 229 predicting the choice of the demonstrator. At the time when the demonstrator's choice was
 230 revealed in the observational learning part, both options were highlighted by a frame to keep
 231 action observation uninformative. Participants received the feedback of their own choice during
 232 the decision phase. They performed the individual learning task in a separate block at the end of
 233 the observational learning sessions. For both tasks, we used ten trials per block, which resulted in
 234 150 trials in total.

235 Behavioral control study. In the behavioral control study, we investigated observational
 236 learning from a non-social agent i.e., from computer-generated choices. The observational
 237 learning task was identical (i.e. instruction, number of sessions, number of trials) to the
 238 observational learning task of the fMRI study described above, except that participants were told
 239 that they observed decisions generated by a computer. The comparison with a computer
 240 demonstrator enabled us to qualify the social observation effect from the fMRI study as outgroup
 241 deficit or ingroup enhancement. In addition to observational learning from the computer,
 242 participants in the control study also performed the individual learning task, for comparison with
 243 the fMRI study. Moreover, we determined participants' political attitude based on their ratings of
 244 the same political initiatives as in the fMRI study. Three participants who showed a right-wing
 245 attitude in these ratings were excluded from the analyses. To keep procedures as similar as
 246 possible to the observational learning task, we did not measure beliefs about the strategy,
 247 programming or reality of the computer generated agent.

248

249 **Statistical analyses of behavioral data**

250 All the behavioral analyses were performed using SPSS version 23.0 (IBM Corp., 2015). For
 251 most of the analyses, we used a repeated measure ANOVA or a paired-t-test because of the
 252 within-subjects design of the study except the κ analysis. When analyzing the κ of ingroup and
 253 outgroup, Wilcoxon sign-ranked tests were used due to the nonparametric nature of the data.

254

255 **Computational modeling**

256 We fitted reinforcement learning models to capture observational learning from actions and
 257 outcomes. Computational modeling was conducted using R (R Core Team, 2013) and a package

258 bbmle. For outcome learning, we assumed that participants track the demonstrator's internal
 259 learning process by simulating outcome prediction errors (sOPE, Eq.1) experienced by the
 260 demonstrator, (Suzuki et al., 2012). These prediction errors correspond to the difference of what
 261 the demonstrator received and the simulated value of the choice made by the demonstrator:

$$\text{Eq 1. } sOPE = R_{t,other_outcome}^{other} - Qa_{t-1,other_outcome}^{other}$$

262 For action learning, we modeled observed action prediction errors (APE, Eq.2) which
 263 relate actually observed choice (1, 0) to the learned probability of observing that choice (Burke et
 264 al., 2010).

$$\text{Eq 2. } APE = A_t^{other} - Pa_{t-1,other_action}^{Subject}$$

265 To model how strongly participants rely on both of these prediction errors to influence
 266 their own choices, we use weights that update the value Q of the chosen option, similar to an
 267 individual learning rate. Specifically, ϖ denotes the degree to which participants incorporate
 268 sOPEs to update Q. The larger the ϖ , the more heavily participants weighted the sOPEs.
 269 Conversely, the action learning weight κ denotes the degree to which participants integrate the
 270 APE into their own choice. We used Equations 3 and 4 for outcome and action learning and
 271 applied a softmax function with a perseverance parameter (range [0 3]) (Wunderlich et al., 2012)
 272 to convert value into action. Both ϖ and κ had a range of [-1 1], allowing for the possibility of
 273 reverse learning from the outgroup.

$$\text{Eq 3. } Qa_{t,other_outcome}^{subject} = Qa_{t-1,other_outcome}^{subject} + \varpi * (R_{t,other_outcome}^{other} - Qa_{t-1,other_outcome}^{other})$$

$$\text{Eq 4. } Pa_{t,other_action}^{Subject} = Pa_{t-1,other_action}^{Subject} + \kappa * (A_t^{other} - Pa_{t-1,other_action}^{Subject})$$

276 We assessed observational learning from the ingroup and outgroup by testing models
 277 with separate or common ϖ and κ for the ingroup and outgroup conditions, resulting in four
 278 models in total. All the models were fitted at the individual level. For model comparison, we
 279 calculated the summed maximum likelihood for all conditions and trials for each participant and
 280 used the Akaike information criterion (AIC) to determine the best model.

281 In addition to the decision for themselves, we modeled the predictions participants made
 282 regarding the choices of the demonstrators, using the most recent sOPE and APE (Eq 5. and Eq
 283 6.) Again, we compared four models with separate or common ϖ and κ for the ingroup and
 284 outgroup conditions to test whether group membership differentially affects prediction learning.
 285 However, we then entered the output of these models into the softmax function to generate
 286 predictions of the decisions of the demonstrator rather than decisions for themselves.

$$\begin{aligned} \text{Eq 5. } Qa_{t+1}^{\text{Subject}}{}^{\text{other_outcome}} &= Qa_t^{\text{Subject}}{}^{\text{other_outcome}} + \varpi * (R_t^{\text{other}}{}^{\text{other_outcome}} - Qa_t^{\text{other}}{}^{\text{other_outcome}}) \\ \text{Eq 6. } Pa_{t+1}^{\text{Subject}}{}^{\text{other_action}} &= Pa_t^{\text{Subject}}{}^{\text{other_action}} + \kappa * (A_t^{\text{other}} - Pa_t^{\text{Subject}}{}^{\text{other_action}}) \end{aligned}$$

289 As an alternative model family, we considered the possibility that action prediction errors
 290 contribute to the learning process in the action-outcome condition as both action and outcome
 291 information are observable in this condition. In order to test whether adding κ -weighted action
 292 prediction errors to the action_outcome condition improves model fit, we examined the model
 293 family described in Equation 7. The weight parameter captures the relative use of sOPE and APE
 294 for updating Q. We again estimated four models varying whether ingroup and outgroup
 295 parameters were common or separate. The average model fit was worse than for the models
 296 without κ in the action_outcome condition (mean AIC: 17.58 vs. 18.56). Moreover, the best-
 297 fitting model of the alternative family explained the data less well than the best-fitting model

without κ in the action_outcome condition. We therefore used the models without κ in the action_outcome condition for further analysis.

$$Eq\ 7. Qa_{t_{other_outcome}}^{Subject} = Qa_{t-1_{other_outcome}}^{Subject} + weight * \varpi * (R_{t_{other_outcome}}^{other} - Qa_{t-1_{other_outcome}}^{other}) + (1 - weight) * \kappa * (A_t^{other} - Pa_{t-1_{other_action}}^{Subject})$$

fMRI data acquisition and analyses

MRI data was acquired with a Philips Achieva 3T whole-body scanner (Philips Medical Systems, Best, The Netherlands) equipped with an 8-channel head coil. For each participant, we collected a T1-weighted whole brain structure image (number of slices: 181, voxel size: 1x1x1mm, field of view: 256x256mm). To measure neural activity we collected T2* weighted whole-brain echo planar images (number of slices: 40, repetition time: 2.36s, voxel size: 3 x 3 x 3mm, field of view: 256x 256mm, echo time: 30ms, flip angle: 90°).

All functional images were distortion corrected, segmented according to the individual T1 structural image, normalized, and smoothed with an 8mm isometric Gaussian kernel (full width at half maximum). Preprocessing and analyses were performed using SPM12 (Wellcome Trust Centre for Neuroimaging). To analyze functional activity, we applied a general linear model with the following regressors: 1) onset of the screen displaying the choice options for the demonstrator, 2) onset of the screen displaying the participant's prediction of demonstrator choice, 3) onset of the screen displaying the choice of the demonstrator, parametrically modulated by 4) the APE (see computational model), 5) onset of the screen displaying the outcome of the demonstrator, parametrically modulated by 6) the sOPE in the action-outcome condition (see computational model), 7) onset of the screen displaying the choice options for the

320 participant, 8) onset of the screen displaying the participant's choice, and 9) onset of the screen
 321 displaying the outcome/masked outcome for the participant. The duration of all events was set to
 322 0. The six head motion regressors and a constant were included as regressors of no interest.

323 We assessed prediction error-related activity in a random effect model with one-sample t-
 324 tests for the contrast images created by the parametric modulators. In order to analyze APE-
 325 related activation independently of demonstrator group, we weighted both ingroup and outgroup
 326 action prediction error regressors with a 1 on the first level and used the resulting contrast images
 327 to perform a one-sample t-test against zero on the second level. The same analysis was
 328 performed to assess sOPE-related activation irrespective of group, using the respective first-level
 329 images from the ingroup and outgroup conditions. We also tested for ingroup vs outgroup
 330 differences in APE and sOPE- related activity at the first-level. Finally, using second-level
 331 correlation, we related the differences in behavioral weights (κ) given to ingroup versus outgroup
 332 action prediction errors to differential neural activity induced from observing ingroup versus
 333 outgroup demonstrator choices. We performed whole brain analyses ($p < 0.05$, family wise
 334 cluster-level whole brain corrected with a cluster inducing voxel-level threshold of $p < 0.001$).

335

336 **Results**

337 Group induction. Prior to scanning, participants rated their own closeness and the closeness of
 338 the future demonstrators in the observational learning task to left- and right-wing parties. A two-
 339 way ANOVA of demonstrator (ingroup/ outgroup) by party (left-wing/ right-wing) revealed
 340 neither a main effect of party ($F(1,28)=1.93.59$, $p=.17$, $\eta^2 = 0.064$) nor demonstrator
 341 ($F(1,28)=1.70$, $p=0.20$, $\eta^2 = 0.057$) but a significant interaction of demonstrator by party
 342 ($F(1,28)=259.59$, $p < 0.001$, $\eta^2 = 0.903$). Participants rated themselves, ($t(28) = 12.64$, $p < 0.001$),

and the future ingroup demonstrator, ($t(28) = 15.31, p < 0.001$), as close to a left-wing party. The future outgroup demonstrator was rated as close to a right-wing party, ($t(28) = 13.70, p < 0.001$) (**Figure 1B**). The difference in closeness ratings to the left-wing party between the participants and the fellow left-wing supporter (ingroup demonstrator) were significantly smaller than the differences in closeness ratings between the participants and the person they perceived as right-wing supporter (outgroup demonstrator), $t(28) = 13.701, p < 0.001$. These results show that the participants perceived one of the demonstrators as a member of their own group (ingroup; defined by left-wing political attitude) and the other demonstrator as a member of a different social group (outgroup; defined by right-wing political attitude).

fMRI study

Behavioral results. In the decision phase of the observational learning task participants made more correct choices (i.e. selected the option associated with higher reward probability) after observing the ingroup compared to the outgroup demonstrator (**Figure. 2A**), group (ingroup, outgroup), $F(1, 28) = 7.839, p = 0.009, \eta^2 = 0.219$. This difference emerged over time, group x trial (one to ten) interaction, $F(9, 252) = 1.938, p = 0.047, \eta^2 = 0.065$. There was no significant difference between the action-only and the action-outcome condition, condition (action_only, action_outcome), $F(1, 28) = 0.988, p = 0.329$, group x condition interaction, $F(1, 28) = 0.820, p = 0.373$, group x condition x trial interaction, $F(9, 252) = 0.925, p = 0.504$. However, separate analyses for each condition revealed that the effect was mainly driven by the action-only condition, showing a significant main effect of group (ingroup, outgroup), $F(1, 28) = 7.421, p = 0.011, \eta^2 = 0.210$ and a significant group x trial interaction, $F(9, 252) = 2.327, p = 0.016, \eta^2 = 0.077$. In the action-outcome condition, we observed a marginally significant main effect of

group, $F(1,28) = 3.323$, $p = 0.079$, $\eta^2 = 0.106$, and no significant group \times trial interaction $p = 0.813$. Compared to individual learning, participants learned less from the outgroup demonstrator, $F(1,28) = 8.168$, $p = 0.008$, $\eta^2 = 0.226$, but similarly well from the ingroup demonstrator, $F(1,28) = 0.174$, $p = 0.714$ (**Figure 2A**). Together, the results show that participants learned less from observing the outgroup compared to the ingroup demonstrator, indicating an outgroup deficit in observational learning, primarily when observing only the actions of others.

Next, we asked if the group difference we observed in choice was mirrored by a similar group difference at the prediction stage, i.e., when participants predicted the upcoming choice of the ingroup and outgroup demonstrator. To test this possibility, we conducted an ANOVA with individual choice predictions and individual choices as dependent variable, and group (ingroup/outgroup) and response type (choice prediction/choice) as independent variables. The results showed a significant main effect of group, $F(1,28) = 7.472$, $p = 0.011$, $\eta^2 = 0.211$, a significant main effect of response type, $F(1,28) = 5.492$, $p = 0.026$, $\eta^2 = 0.164$, and a significant group \times response type interaction, $F(1,28) = 5.35$, $p = 0.028$, $\eta^2 = 0.160$ (**Figure 2B**). Clarifying this interaction effect, post-hoc pairwise comparisons showed that participants predicted the choices of the ingroup and the outgroup demonstrators equally well, $t(28) = 1.408$, $M_{\text{difference}} = 0.027$, $SE = 0.019$, $p = 0.171$, but showed significantly fewer correct choices in the outgroup compared to the ingroup condition, $t(28) = 2.800$, $M_{\text{difference}} = 0.112$, $SE = 0.040$, $p = 0.009$. Thus, participants learned to predict ingroup and outgroup choices similarly well, but used the learned information to a lesser degree when observing the outgroup compared to the ingroup demonstrator.

Behavioral control study

389 We conducted a behavioral control experiment to clarify whether the observed ingroup vs
 390 outgroup difference in observational learning reflects increased learning from the ingroup or
 391 reduced learning from the outgroup. The control experiment was identical to the main
 392 experiment, except that participants observed choices from a computer, i.e., a non-human agent.
 393 In case of increased learning from the ingroup demonstrator, the number of correct choices
 394 should be significantly higher in the ingroup demonstrator condition compared to the computer
 395 condition. Conversely, in case of decreased learning from the outgroup demonstrator, the number
 396 of correct choices should be significantly lower in the outgroup demonstrator condition
 397 compared to the computer condition. To test this issue, we performed two repeated measures
 398 ANOVAs with trials and condition (action_only/action_outcome) as within subject variables and
 399 demonstrator (ingroup/outgroup, computer) as a between subject variable. The results showed
 400 significantly fewer correct choices after observing the outgroup demonstrator, compared to the
 401 computer, demonstrator, $F(1,57) = 14.343$, $p = 0.0003$, $\eta^2 = 0.201$, mean difference between the
 402 outgroup and the computer conditions: $M_{\text{combined}} = 0.142$, $SE = 0.029$; $M_{\text{action_outcome}} = 0.135$; SE
 403 $= 0.035$, $M_{\text{action_only}} = 0.205$, $SE = 0.038$ (see **Figure 2C**). There were no other significant effects,
 404 condition (action_only / action_outcome), $F(1,57) = 0.464$, $p = 0.499$, $\eta^2 = 0.008$, demonstrator
 405 x condition interaction, $F(1,57) = 1.66$, $p = 0.203$, $\eta^2 = 0.008$. In contrast to the difference
 406 between outgroup and computer demonstrator, there were no significant differences in the
 407 number of correct choices between the ingroup and the computer condition, demonstrator
 408 (ingroup / computer), $F(1,57) = 1.639$, $p = 0.208$, $\eta^2 = 0.028$, condition (action_only /
 409 action_outcome), $F(1,57) = 0.090$, $p = 0.765$, $\eta^2 = 0.002$, demonstrator x condition interaction,
 410 $F(1,57) = 0.174$, $p = 0.678$, $\eta^2 = 0.003$, mean difference between ingroup and computer conditions:
 411 $M_{\text{combined}} = 0.057$, $SE = 0.029$, $M_{\text{action_outcome}} = 0.048$, $SE = 0.031$, $M_{\text{action_only}} = 0.067$, $SE =$

0.036 (see **Figure 2C**). Individual learning in the fMRI study and control study were not different, ($F(1,58) = 1.031$, $p = .314$, $\eta^2 = .016$). These results indicate that the observed group difference in observational learning reflects an outgroup deficit, rather than enhanced learning from the ingroup.

Computational modeling

We hypothesized that an outgroup deficit in observational learning might be driven by differences in outcome-related learning, differences in action-related learning, or differences in both learning mechanisms. To test these hypotheses, we fitted reinforcement learning models to choice behavior when participants observed outcomes and actions from ingroup or outgroup demonstrators. Model comparisons showed that behavior was best characterized by a model that used a common learning weight (ϖ) for ingroup and outgroup outcome prediction errors (OPEs), but separate learning weights (κ) for ingroup and outgroup action prediction errors (APEs; **Figure 3A**). The κ -ingroup weight was larger than the κ -outgroup weight, Wilcoxon Rank-sum test, $z = 3.13$, $p = .002$ (**Figure 3B**). Thus, model comparison results support the notion that the differential observational learning effect is mainly due to reduced action-based learning from the outgroup.

Using a similar approach and models, we performed model comparisons also for participants' predictions of the demonstrators' choices. The best model for prediction behavior was the model with common learning weight (ϖ) for ingroup and outgroup outcome prediction errors (OPEs) and learning weight (κ) for ingroup and outgroup action prediction errors (**Figure 3C**). Thus, our participants learned to predict the choices of ingroup and outgroup demonstrators similarly well (but used the acquired information differentially for their own choices).

435

436 **fMRI**

437 *Replication of previous results: Observational learning irrespective of group membership.* First,
 438 we investigated if our neural results replicate the findings of previous studies that investigated
 439 observational learning irrespective of group membership, showing activation in MPFC
 440 associated with outcome prediction errors and activation in DLPFC related to action prediction
 441 errors (Burke et al., 2010; Suzuki et al., 2012). To do so, we conducted two separate parametric
 442 regression analyses that regressed the participants' trial-by-trial model estimates of ingroup and
 443 outgroup outcome-prediction errors and action prediction errors against their neural activity
 444 during the observation of outcomes or actions, respectively.

445 Group independent outcome-related learning activated a network of brain regions,
 446 including the dorsomedial prefrontal cortex (DMPFC), bilateral insula, caudate and midbrain
 447 (**Table 1, Figure 4A and B**). The neural response in these regions increased with decreasing
 448 outcome-prediction errors, indicating that observed outcomes eliciting smaller outcome-
 449 prediction errors resulted in stronger activity. At the applied threshold, there was no region
 450 where neural response increased with increasing outcome-prediction errors. Notably, the
 451 DMPFC findings co-localize with previously unreported findings of inverse outcome prediction
 452 error coding in the study of Burke and colleagues (2010; **Figure 4A**).

453 Conversely, learning from observing ingroup and outgroup actions activated the anterior
 454 IFG and parietal regions. (**Table 2, Figure 5A and B**). The IFG findings co-localize with those
 455 of our previous report on action prediction error coding in that area (Burke et al., 2010, **Figure**
 456 **5A**).

457

458 *Group differences in observational learning mechanisms.* Second, we investigated the brain
 459 regions that are differentially involved in learning from ingroup and outgroup outcomes and
 460 actions. To do so, we contrasted the neural responses related to participants' trial-by-trial model
 461 estimates of ingroup outcome-prediction errors and action prediction errors with the neural
 462 responses related to their trial-by-trial model estimates of outgroup outcome-prediction errors
 463 and action prediction errors. The results revealed no significant differences at the applied
 464 threshold, suggesting that participants activated similar neural circuitries while learning from the
 465 observation of ingroup and outgroup outcomes and actions.

466
 467 *Group differences in the weight assigned to observational learning.* Third, we tested our
 468 assumption that the stronger weight assigned to action prediction errors in the ingroup compared
 469 to the outgroup condition (**Fig. 3B**) is related to neural activation of the DLPFC/IFG. Using a
 470 second-level regression analysis, we regressed the behavioral contrast between ingroup and
 471 outgroup action learning weights κ against the neural contrast between the observation of
 472 ingroup and outgroup choices, i.e., the time when observational action prediction error can be
 473 computed. The results revealed only one significant whole brain corrected result, in the left IFG
 474 (MNI_{xyz} : -34, 0, 28, Z_{stats} = 4.07, P_{FWE} whole-brain corrected = 0.041) (**Figure 6A and B**). Thus, left IFG
 475 activity reflected the impact of action prediction errors on behavior, which was reduced when
 476 participants observed outgroup actions compared to when they observed ingroup actions. This
 477 activity localized in the posterior part of the IFG region that was identified as the key region of
 478 the action prediction error learning network.

479 480 Discussion

481 In this study, we investigated whether observational learning is shaped by the important social
482 factor group membership. We report novel evidence that participants learn similarly well from
483 observing ingroup and outgroup outcomes, but learn less well from observing outgroup actions.
484 The observed deficit in learning from outgroup actions provides a plausible source of individuals'
485 difficulties to learn from outgroup members, a phenomenon that has been described in previous
486 studies (Buttelmann et al., 2013; Golkar et al., 2015; Golkar and Olsson, 2017), but so far not
487 explained. Our neural results not only converge with the findings of previous observational
488 learning studies in lateral prefrontal cortex (Burke et al., 2010), but also showed that IFG
489 differentially encodes learning from observing ingroup vs outgroup actions.

490 In more detail, our behavioral findings revealed that participants made fewer correct
491 choices after observing an outgroup demonstrator, compared to an ingroup demonstrator (Figure
492 2A) or to a computer demonstrator (Figure 2C). These outgroup deficits in observational learning
493 occurred although participants observed comparable choice behavior in the ingroup and the
494 outgroup condition and were able to predict ingroup and outgroup choices equally well. The
495 finding of reduced learning from observing an outgroup compared to an ingroup individual is in
496 line with previous behavioral evidence (Buttelmann et al., 2013; Golkar et al., 2015).

497 Extending these previous studies, we used computational learning models to specify the
498 source of the outgroup deficit in observational learning. Given that observational learning is
499 based on learning from observed outcomes (Burke et al., 2010; Suzuki et al., 2012; Kumaran et
500 al., 2015) and observed actions (Burke et al., 2010; Suzuki et al., 2012), we hypothesized that
501 outgroup deficits in observational learning might occur because individuals rely more on
502 observed ingroup compared to outgroup outcomes, reflected by a stronger weight for ingroup
503 compared to outgroup outcome prediction errors. Alternatively, we assumed that outgroup

504 deficits in observational learning might arise because participants rely more on observed ingroup
 505 compared to outgroup actions, reflected by a stronger weight for ingroup compared to outgroup
 506 action prediction errors. Our computational modeling results showed that the behavioral
 507 outgroup deficit in observational learning were best explained by a model in which participants
 508 put less weight on action prediction errors elicited by an outgroup individual than on action
 509 prediction errors elicited by an ingroup individual (Figure 3) and put similar weight on outcome
 510 prediction errors from the two groups.

511 The computational modeling results converge with the behavioral findings that showed a
 512 clear ingroup vs outgroup difference in the condition in which participants could only learn from
 513 actions, i.e., the action_only condition, and a marginal main effect of group ($F(1,28) = 3.32$, $p =$
 514 0.079) in the condition in which participants could also learn from outcomes, i.e., the
 515 action_outcome condition. Presumably, the group effect in the action_outcome condition did not
 516 reach significance because the behavioral results reflect a mix of outcome-related, i.e., unbiased,
 517 and action-related, i.e., biased, learning. As computational models establish a relation between
 518 components of the phenomenon being modeled and the components of the model (Stafford,
 519 2009), they can be more sensitive to the latent processes that might drive modulations in
 520 behavior than statistical analyses comparing behavioral outcomes alone (Stafford et al., 2020). In
 521 line with this notion, using computational modeling allowed us to disentangle the mixture of
 522 outcome- and action-based learning, and to specify the effect of group membership on the
 523 different subcomponents of observational learning.

524 Our neural results revealed that action- and outcome prediction errors elicited by observing
 525 the ingroup and the outgroup demonstrator are processed by similar neural circuitries that
 526 replicate previous findings. In more detail, learning from outcome prediction errors was

527 associated with activation in the DMPFC, the insula, the caudate and the midbrain, i.e., regions
 528 that have been implicated in outcome-related learning in previous neuroscience studies (Liu et al.,
 529 2011; Sescousse et al., 2013). Conversely, learning from action prediction errors was linked to
 530 neural responses in the anterior portion of the inferior frontal gyrus (IFG) and the parietal cortex,
 531 again in line with previous evidence (Burke et al., 2010; Suzuki et al., 2012). Key regions of the
 532 outcome-learning and action-learning networks that we obtained in the current studies showed
 533 considerable overlap with the respective neural circuitries observed in an independent previous
 534 study (Burke et al., 2010) that investigated observational learning independent of group
 535 membership (Figure 4).

536 Interestingly, although learning from ingroup and outgroup prediction errors activated
 537 similar neural networks, participants put stronger weight on the use of ingroup than outgroup
 538 action prediction error when they made decisions for themselves, which was reflected by
 539 stronger activation in the IFG (Figure 6). The IFG is involved in action-observation and imitation
 540 processes (Caspers et al., 2010) and forms part of the mirror neuron system (Molenberghs et al.,
 541 2012). Moreover, there is evidence that the activity of this area is modulated by group
 542 membership. For example, greater IFG activity was found when participants evaluated an
 543 ingroup member based on detailed personal information as compared to an outgroup member
 544 (Freeman et al., 2010). Another recent neuroimaging study revealed stronger activation in a
 545 mirror neuron network including left IFG when participants observed facial emotions of ingroup
 546 individuals compared to outgroup individuals (Krautheim et al., 2019). In line with this previous
 547 evidence, our results show that the processing of perceived actions in the IFG is modulated by
 548 group membership of the demonstrator during observational learning. Extending these previous
 549 findings, our results indicate that the IFG activity selects action-related information based on

550 social information (here group membership), and thus forms a plausible neural basis for biases in
 551 observational learning. However, the putative link with the mirror system will need to be tested
 552 formally.

553 It is worth noting the limitations of our study. First, we used exclusively political attitude
 554 to manipulate group membership. Given that previous research described outgroup learning
 555 deficits with group membership based on language or race (Buttelmann et al., 2013; Golkar et al.,
 556 2015), it is unlikely that outgroup learning deficits are limited to the political domain. Moreover,
 557 it is well-established that differences in political attitude foster social categorization, i.e., the
 558 formation of social ingroups and outgroups (Caruso et al., 2009; Rand et al., 2009; Young et al.,
 559 2014). There is even evidence that differences in political attitude can override social
 560 categorization based on race (Losin and Woo, 2015). In line with this previous evidence, the
 561 group induction based on political attitude in our study resulted in a salient group membership
 562 manipulation (Figure 1B), a conclusion further supported by our findings of significant
 563 behavioral and neural differences between the ingroup and the outgroup condition. That said,
 564 future research may want to investigate observational learning with a different group
 565 manipulation. Second, we studied left-wing participants only. Although targeting only one group
 566 (e.g. white participants in a study on race) to investigate ingroup-outgroup behavior is common
 567 in the literature (Golkar et al., 2015; Hein et al., 2016), future research may also want to study
 568 right-wing individuals or other political groups to generalize our findings. Third, we used a
 569 relatively small sample size and recent studies (Bossier et al., 2020; Marek et al., 2020)
 570 recommend larger sample sizes to ensure replication of fMRI findings than current practice
 571 (Yeung, 2018). It is therefore noteworthy that we replicate previous findings (Burke et al., 2010)
 572 on group-independent observational learning. Still, particularly our correlation findings should

573 be re-assessed with a larger sample in the future. Fourth, we did not specify the outgroup
 574 attributes that drive or shape the group differences in observational learning mechanisms
 575 revealed in our study. Future research may want to investigate whether the observed difference
 576 in observational learning arose because participants are less likely to trust outgroup actions
 577 without disambiguating feedback about the correctness of the outgroup demonstrator's choice
 578 (i.e., the outgroup outcome). Another factor that might play a role is the extent to which
 579 participants dislike the outgroup. In line with the findings of other studies (Golkar et al., 2015;
 580 Hein et al., 2016) it is conceivable that the individual impressions and/or emotions towards the
 581 respective outgroup might modulate the outgroup-related observational learning deficits
 582 observed in our study, an assumption that should be investigated in future studies.

583 In conclusion, the current results reveal that outgroup deficits in observational learning
 584 mainly reflect decreased learning from observed actions. Our findings suggest that the IFG
 585 differentially weighs ingroup and outgroup action prediction errors and provide a
 586 neurocomputational mechanism for outgroup deficits in observational action learning.

587 **Data and code availability**

588 The behavioral data is available online, <https://osf.io/savw4/>. Also, the neuroimaging results in
 589 this study can be found in <https://neurovault.org/collections/VUOYOUFT/>. Code to implement
 590 computational model is available from the corresponding author upon reasonable request.

591 **References**

- 592 Bossier H et al. (2020) The empirical replicability of task-based fMRI as a function of sample
 593 size. *Neuroimage* 212:1–12.
 594 Burke CJ, Tobler PN, Baddeley M, Schultz W (2010) Neural mechanisms of observational
 595 learning. *Proc Natl Acad Sci* 107:14431–14436 Available at:
 596 <http://www.pnas.org/cgi/doi/10.1073/pnas.1003111107>.

- Buttelmann D, Zmyj N, Daum M, Carpenter M (2013) Selective Imitation of In-Group Over Out-Group Members in 14-Month-Old Infants. *Child Dev* 84:422–428.
- Caruso EM, Mead NL, Balcetis E (2009) Political partisanship influences perception of biracial candidates' skin tone. *Proc Natl Acad Sci U S A* 106:20168–20173.
- Caspers S, Zilles K, Laird AR, Eickhoff SB (2010) ALE meta-analysis of action observation and imitation in the human brain. *Neuroimage* 50:1148–1167 Available at: <http://dx.doi.org/10.1016/j.neuroimage.2009.12.112>.
- Charpentier CJ, Iigaya K, O'Doherty JP (2020) A Neuro-computational Account of Arbitration between Choice Imitation and Goal Emulation during Human Observational Learning. *Neuron*:1–13 Available at: <https://doi.org/10.1016/j.neuron.2020.02.028>.
- Freeman JB, Schiller D, Rule NO, Ambady N (2010) The neural origins of superficial and individuated judgments about ingroup and outgroup members. *Hum Brain Mapp* 31:150–159.
- Golkar A, Castro V, Olsson A (2015) Social learning of fear and safety is determined by the demonstrator's racial group. *Biol Lett* 11:20140817 Available at: <http://rsbl.royalsocietypublishing.org/content/11/1/20140817>.
- Golkar A, Olsson A (2017) The interplay of social group biases in social threat learning. *Sci Rep* 7:1–5.
- Hein G, Engelmann JB, Vollberg MC, Tobler PN (2016) How learning shapes the empathic brain. *Proc Natl Acad Sci* 113:80–85 Available at: <http://www.pnas.org/lookup/doi/10.1073/pnas.1514539112>.
- Howard LH, Henderson AME, Carrazza C, Woodward AL (2015) Infants' and Young Children's Imitation of Linguistic In-Group and Out-Group Informants. *Child Dev* 86:259–275.
- IBM Corp. (2015) IBM SPSS Statistics for Windows.
- Krautheim JT, Dannlowski U, Steines M, Neziroğlu G, Acosta H, Sommer J, Straube B, Kircher T (2019) Intergroup empathy: Enhanced neural resonance for ingroup facial emotion in a shared neural production-perception network. *Neuroimage* 194:182–190.
- Kumaran D, Warren DE, Tranel D (2015) Damage to the Ventromedial Prefrontal Cortex Impairs Learning from Observed Outcomes. *Cereb Cortex* 25:4504–4518 Available at: <https://academic.oup.com/cercor/article-lookup/doi/10.1093/cercor/bhv080>.
- Liu X, Hairston J, Schrier M, Fan J (2011) Common and distinct networks underlying reward valence and processing stages: A meta-analysis of functional neuroimaging studies. *Neurosci Biobehav Rev* 35:1219–1236 Available at: <http://dx.doi.org/10.1016/j.neubiorev.2010.12.012>.
- Losin EAR, Iacoboni M, Martin A, Cross KA, Dapretto M (2012) Race modulates neural activity during imitation. *Neuroimage* 59:3594–3603.
- Losin EAR, Woo C (2015) Brain and psychological mediators of imitation : sociocultural versus physical traits. *Cult Brain*:93–111 Available at: <http://dx.doi.org/10.1007/s40167-015-0029-9>.
- Marek S et al. (2020) Towards Reproducible Brain-Wide Association Studies. *bioRxiv* 11:15–18 Available at: <https://doi.org/10.1101/2020.08.21.257758>.
- Molenberghs P, Cunnington R, Mattingley JB (2012) Brain regions with mirror properties: A meta-analysis of 125 human fMRI studies. *Neurosci Biobehav Rev* 36:341–349 Available at: <http://dx.doi.org/10.1016/j.neubiorev.2011.07.004>.

- 642 R Core Team (2013) R: A language and environment for statistical computing. Available at:
643 <http://www.r-project.org/>.
- 644 Rand DG, Pfeiffer T, Dreber A, Sheketoff RW, Wernerfelt NC, Benkler Y (2009) Dynamic
645 remodeling of in-group bias during the 2008 residential election. *Proc Natl Acad Sci U S A*
646 106:6187–6191.
- 647 Sescousse G, Caldú X, Segura B, Dreher JC (2013) Processing of primary and secondary
648 rewards: A quantitative meta-analysis and review of human functional neuroimaging studies.
649 *Neurosci Biobehav Rev* 37:681–696 Available at:
650 <http://dx.doi.org/10.1016/j.neubiorev.2013.02.002>.
- 651 Suzuki S, Harasawa N, Ueno K, Gardner JL, Ichinohe N, Haruno M, Cheng K, Nakahara H
652 (2012) Learning to Simulate Others' Decisions. *Neuron* 74:1125–1137 Available at:
653 <http://dx.doi.org/10.1016/j.neuron.2012.04.030>.
- 654 Wunderlich K, Smittenaar P, Dolan RJ (2012) Dopamine Enhances Model-Based over Model-
655 Free Choice Behavior. *Neuron* 75:418–424.
- 656 Yeung AWK (2018) An updated survey on statistical thresholding and sample size of fMRI
657 studies. *Front Hum Neurosci* 12:1–7.
- 658 Young AI, Ratner KG, Fazio RH (2014) Political Attitudes Bias the Mental Representation of a
659 Presidential Candidate's Face. *Psychol Sci* 25:503–510.
- 660

661

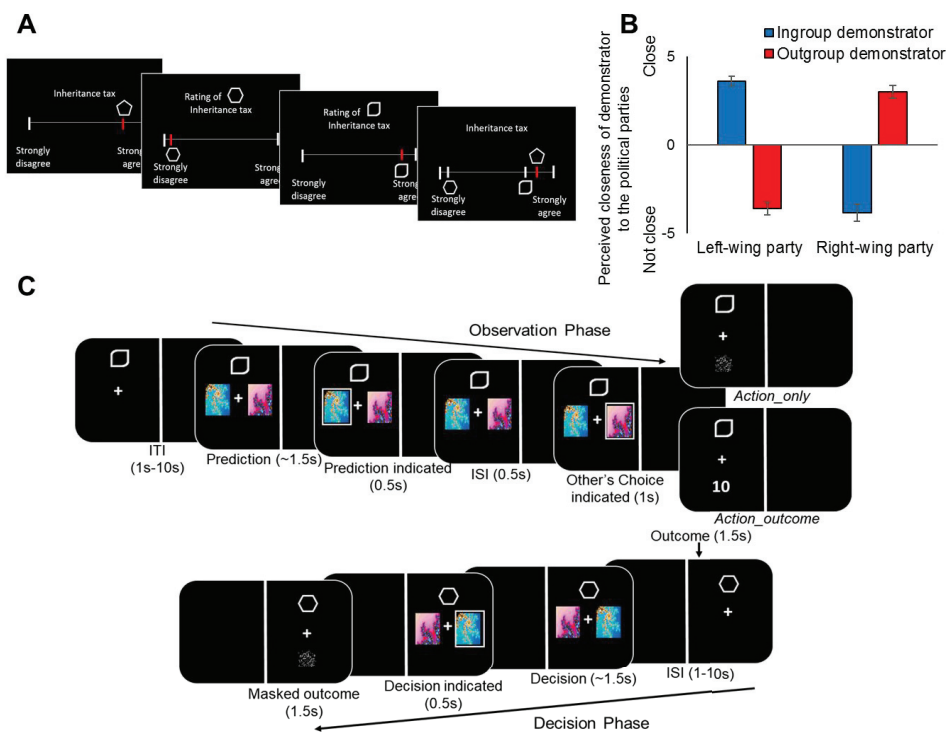
Table 1. Summary of brain regions correlating with -sOPE as a parametric modulator in both ingroup and outgroup conditions, whole brain cluster-level FWE- corrected $p < 0.05$.

Brain region	Coordinates			T value	Z value	Voxels in cluster
	X	Y	Z			
DMPFC	-2	8	54	5.06	4.23	985
Caudate (L)	-16	10	6	4.30	3.73	283
Insula (L)	-42	16	4	4.73	4.02	176
Insula (R)	38	6	-8	4.62	3.95	263
Midbrain	-10	-36	-34	4.61	3.94	262
Precuneus	-26	-54	46	4.61	3.94	247

Table 2. Summary of brain regions correlating with APE as a parametric modulator in both ingroup and outgroup conditions, whole brain cluster-level FWE- corrected $p < 0.05$.

Brain region	Coordinates			T value	Z value	Voxels in cluster
	X	Y	Z			
Inferior Frontal gyrus (L)	-46	8	34	5.12	4.27	951
Inferior Frontal gyrus (R)	42	6	34	4.94	4.16	799
Precuneus/Inferior parietal lobe (R)	34	-68	40	4.88	4.11	467
Inferior parietal lobe (L)	-38	-48	44	4.79	4.06	654
Cerebellum /occipital lobe (R)	30	-62	0	5.70	4.61	1794
Cerebellum /occipital lobe (L)	-34	-74	-20	5.47	4.47	1463

673 **Figure captions**

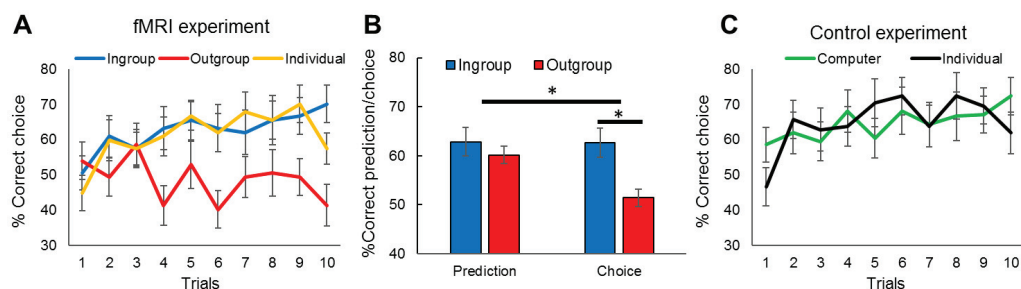


674

675 **Figure. 1. Group induction, manipulation check, and observational learning task. A)** Prior
 676 to scanning, participants rated current political issues in Switzerland and viewed ratings of two
 677 different persons (future demonstrators) regarding the same issues. **B)** Participants' ratings of the
 678 ingroup and outgroup demonstrator's closeness towards a left-wing and right-wing party, based
 679 on how the demonstrators rated political initiatives (as in the example shown in Figure 1A). The
 680 participants rated the ingroup demonstrator as similar to themselves, i.e., close to a left-wing
 681 party and the outgroup demonstrator as dissimilar to themselves, i.e., close to a right-wing
 682 party. Blue = perceived closeness of the ingroup demonstrator towards a left-wing (left panel) or right-
 683 wing (right panel) party; red = perceived closeness of the outgroup demonstrator towards a left-
 684 wing (left panel) or right-wing (right panel) party. Error bars indicate standard errors of the
 685 mean. **C)** In the main part of the study, participants observed the ingroup or outgroup
 686 demonstrator choosing between two fractal images. First, the demonstrator was indicated by one
 687 of two abstract symbols (counterbalanced across participants). Then, participants had 1.5 s to

688 predict which option the demonstrator would choose. After response registration (0.5s) and a
 689 short interval (0.5s), the action of the demonstrator was shown for 1s. Depending on the
 690 condition, participants observed the outcome of the demonstrator (action_outcome condition) or
 691 a pixel-matched scrambled image (action_only condition) for 1.5 s. Next, the demonstrator's
 692 symbol was replaced by the participant's symbol, and the participants chose between the same
 693 options.

694
 695
 696
 697
 698



699

700 **Figure 2. Behavioral results. A)** Trial-wise percentage of correct choices in the three
 701 experimental conditions (ingroup, outgroup and individual learning) of the fMRI study **B)**
 702 Average percentage of correct choices and predictions. Group membership of the demonstrator
 703 affected correct choice but not prediction. **C)** Trial-wise percentage of correct choices in the two
 704 conditions (computer and individual learning) of the control experiment. Given that there were
 705 no significant differences between the action_only and the action_outcome conditions, the results
 706 are pooled over these two conditions in (A) and (C). Error bars indicate standard errors of the
 707 mean.

708
 709

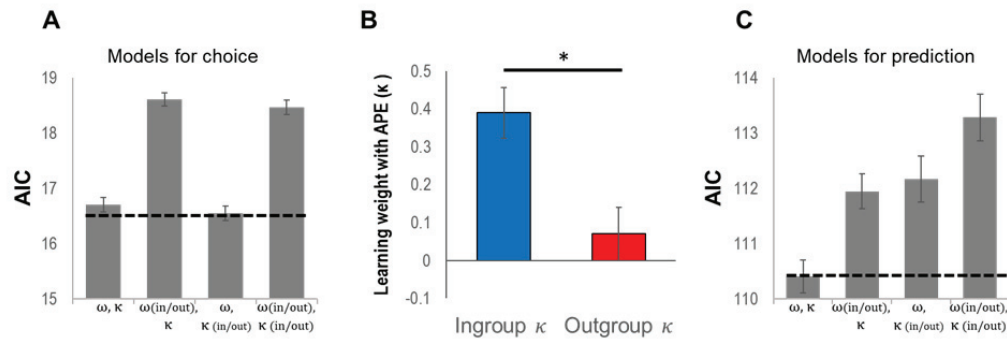


Figure 3. Computational modeling results. **A)** Model comparison based on the Akaike information criterion (AIC) favored the model with a common outcome learning weight ω and separate action learning weights κ for ingroup and outgroup demonstrators when participants made decision for themselves. **B)** Ingroup vs outgroup difference in action learning weight (κ). The weight given to ingroup action prediction errors (κ) was larger than the one given to outgroup action prediction errors in behavior. Error bars indicate standard errors of the mean. **C)** Model comparison for predictions of demonstrators' choices. AIC favored the model with a common outcome learning weight ω and common action learning weights κ for ingroup and outgroup demonstrators when participants predicted the decisions of the demonstrators.

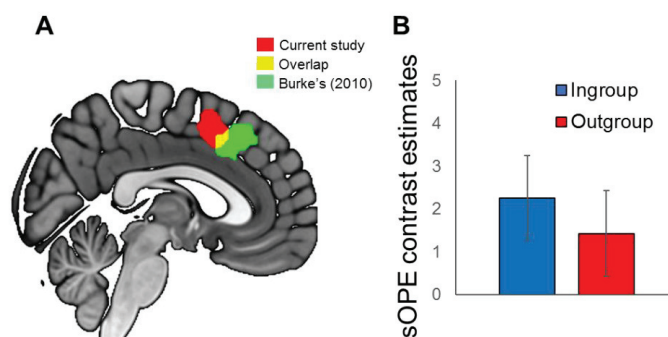


Figure 4. Outcome prediction error coding irrespective of group. **A)** Activation in dorsomedial prefrontal cortex (DMPFC) correlating with inverse simulated outcome prediction errors (sOPE) in both ingroup and outgroup conditions (red). This region overlapped (yellow) with the DMPFC region that correlated with inverse simulated outcome prediction errors in an independent previous study (green; unpublished data from Burke et al., 2010). For illustration purposes, results are displayed at p uncorrected < 0.001 (see **Table 1** for details and whole brain results). **B)** Bar plot of the DMPFC region shown in red (A), illustrating activity correlating with inverse simulated outcome prediction error for both ingroup and outgroup demonstrators. Error bars indicate standard errors of the mean.

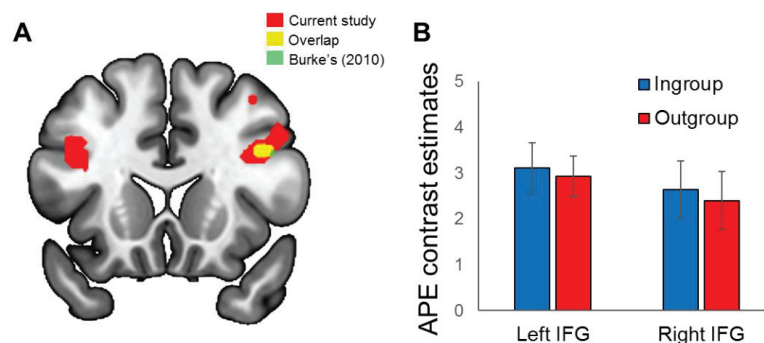
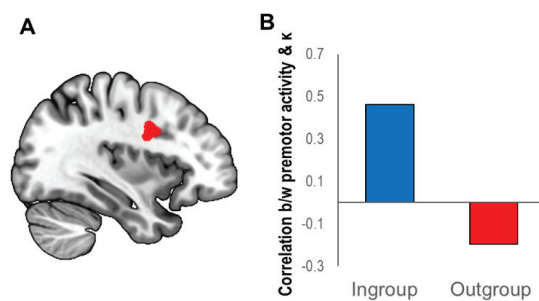


Figure 5. Action prediction error coding irrespective of group. **A)** Prefrontal regions in left and right IFG correlated with action prediction error (APE) as a parametric modulator regardless of group (red). The region in right IFG overlapped with the DLPFC regions that correlated with action prediction errors in a previous independent study (Burke et al., 2010; yellow / green). For illustration purposes, results are displayed at p uncorrected < 0.001 (see **Table 2** for details and whole brain results). **B)** Bar plots illustrating the relationship between the activity in left and right IFG regions shown in (A) and ingroup and outgroup action prediction errors. Error bars indicate standard errors of the mean.

747



748

749 **Figure 6. Group differences in neural activity correlating with behavioral effects.** A) The
 750 ingroup vs. outgroup difference in κ correlated with ingroup vs. outgroup differences in inferior
 751 frontal gyrus (IFG, family wise cluster-level whole brain correction, $p = 0.041$, with an
 752 uncorrected voxel-level (i.e., cluster-inducing) threshold of $p < 0.001$). For illustration purposes,
 753 results are displayed at p uncorrected < 0.001 . B) Correlation between individual IFG activity
 754 (extracted from the cluster) and κ .

